

# Resolving Ambiguity in Embodied Instructions via Semantic Valency Conflict

Anonymous ACL submission

## Abstract

Natural language instructions in human-robot interaction often contain subtle ambiguities that hinder reliable interpretation. These ambiguities arise when a single instruction can be interpreted in multiple ways, assigning conflicting semantic roles to objects, tools, or participants, potentially leading to execution failures. To address this, we propose Semantic Valency Conflict (SVC), a cognitively inspired, logit-free method for detecting ambiguity in robot-directed instructions. SVC identifies divergences in role assignments across alternative interpretations of a predicate, using large language models (LLMs) to infer context-sensitive semantic frames. Our method is model-agnostic and compatible with both open- and closed-source LLMs. SVC produces clear, structured outputs that highlight which parts of the instruction are ambiguous and indicate which predicate and its associated arguments lead to multiple or conflicting interpretations. We evaluate SVC on two datasets, AmbiK and Introspective Planning, and demonstrate that it shows strong and consistent performance in detecting subtle ambiguities in natural language instructions given to robots in safety, unambiguous, and preference-based scenarios.

## 1 Introduction

Natural language is inherently polysemous, making it challenging for large language models (LLMs) to accurately follow instructions (Heo et al., 2025), interpret textual descriptions (Singhal et al., 2024; Chuganskaya et al., 2023), and perform planning tasks (Hazra et al., 2024; Grigorev et al., 2024, 2025). Depending on the context, multiple valid actions may exist in a given situation, making it difficult to select the most appropriate one. These challenges are associated with two types of uncertainty: epistemic, arising from lack of knowledge or insufficient data, and aleatoric, caused by random and unpredictable variations in the environ-

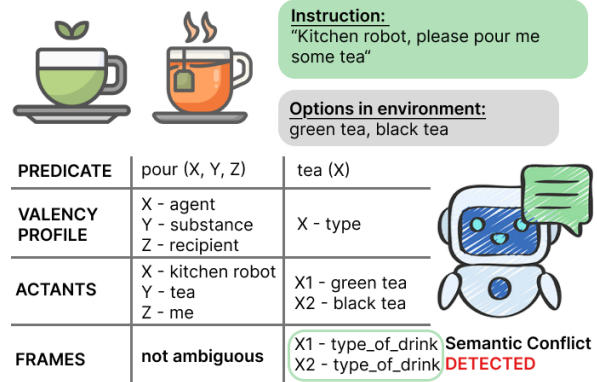


Figure 1: The robot receives the instruction “Pour me some tea” in a setting with green and black tea. The method detects ambiguity by identifying competing interpretations of the predicate “tea”, where “green” and “black” act as conflicting actants within the same action frame, revealing a preference-related variability.

ment (Shorinwa et al., 2024). Both types of uncertainty complicate action selection by introducing ambiguity in outcome evaluation and adaptation to new situations.

In robotics, one of the key challenges is following instructions given in natural language. Modern approaches to solving this problem primarily rely on LLMs, which, unlike heuristic methods, can flexibly interpret complex commands (Huang et al., 2022; Ahn et al., 2022; Kovalev and Panov, 2022; Sarkisyan et al., 2023; Onishchenko et al., 2025). However, user instructions can often be ambiguous due to factors such as the use of synonyms, metaphorical expressions, abbreviations, or environmental complexity where the same object may be represented in various forms (Fig. 1). This ambiguity, alongside known risks like LLM hallucinations, can lead to mission failures or safety hazards (Zhang et al., 2025). To mitigate these risks, it is essential to detect ambiguity promptly and clarify instructions. Implementing ambiguity detection plays a vital role in this process, ensuring the safe

065 integration of LLMs into intelligent agents (Firoozi  
066 et al., 2023). Accurate interpretation and execution  
067 of instructions are therefore crucial for building  
068 reliable and safe robotic systems.

069 Common techniques for ambiguity resolution in-  
070 clude generating multiple candidate interpretations  
071 ranked by contextual relevance and interactive clar-  
072 ification (Ren et al., 2023). Some approaches also  
073 incorporate external knowledge bases to improve  
074 accuracy, but this reduces the system’s autonomy  
075 and adaptability in novel environments (Liang et al.,  
076 2024). Although recent methods apply techniques  
077 like Conformal Prediction (CP) such as Su et al.  
078 (2024a) and affordance estimation Jr. and Manocha  
079 (2024) to improve ambiguity detection, they still  
080 struggle to effectively capture subtle ambiguities  
081 inherent in natural language instructions.

082 To address ambiguity in natural language instruc-  
083 tions caused by lexical polysemy and underspeci-  
084 fied argument structures, we propose the **Semantic  
085 Valency Conflict (SVC)** method. The key idea is  
086 that a single word (or lexeme) can trigger multiple,  
087 incompatible interpretations known as semantic  
088 frames, each expecting a different set of roles and  
089 participants to make sense in context (also called  
090 a valency profile, see Fig. 1). By modeling ambi-  
091 guity as a conflict between these profiles under a  
092 given environment (set of objects, properties, and  
093 relations present in the external context), SVC de-  
094 termines whether an instruction supports multiple  
095 mutually exclusive interpretations.

096 The SVC method identifies lexical units whose  
097 frame-induced valency profiles are incompatible  
098 within the current environment. This is achieved  
099 through predicate identification, dynamic frame  
100 induction and contextual role alignment.

101 Unlike prior approaches, SVC detects fine-  
102 grained semantic conflicts without logit access or  
103 static knowledge bases, enabling interpretability  
104 and adaptability in both white- and black-box set-  
105 tings. Fine-grained ambiguity refers to subtle se-  
106 mantic inconsistencies arising at a detailed level  
107 of meaning, such as nuanced conflicts between se-  
108 mantic roles, argument structures, or contextual  
109 interpretations.

110 We evaluated SVC method ability to detect am-  
111 biguity on two fully textual datasets: IntroPlan  
112 Mobile Manipulation (Liang et al., 2024) and Am-  
113 biK (Ivanova et al., 2025). Compared to logit-based  
114 and heuristic baselines, which either fail to detect  
115 ambiguity or ignore uncertain cases entirely, SVC  
116 demonstrates superior performance in identifying

117 semantically ambiguous instructions, particularly  
118 in user-preference contexts.

119 **Statement of contributions.** In this work, we  
120 propose Semantic Valency Conflict, a new method  
121 for ambiguity detection grounded in frame-based  
122 valency analysis. We develop logit-free architec-  
123 ture that is compatible with both white-box and  
124 black-box language models. Unlike prior work,  
125 our method with LLM dynamically induces cog-  
126 nitive frames without relying on static knowledge  
127 bases or plans, making it adaptable. Finally, SVC  
128 operates at the cognitive layer of instruction under-  
129 standing, ensuring semantic reliability before plan-  
130 ning. Its structured outputs can directly inform or  
131 condition the planner by providing disambiguated  
132 action frames, reducing misinterpretations that lead  
133 to execution errors.

## 134 2 Related Works

135 LLMs exhibit uncertainty due to external and inter-  
136 nal factors. External factors, like data noise, ambi-  
137 guity, and lack of contextual information can some-  
138 times be mitigated with clarifying questions (Zhang  
139 and Choi, 2025). Internal factors, such as model  
140 architecture, training limitations, and probabilistic  
141 text generation can be reduced through model re-  
142 finement or by increasing the number and diversity  
143 of in-context examples (Wang et al., 2025).

144 Various methods are used for resolving ambi-  
145 guity, including token-based, self-verbalized,  
146 semantic-similarity, and mechanistic interpretabil-  
147 ity methods (Shorinwa et al., 2024). Model deci-  
148 sion analysis using methods such as directed en-  
149 tailment graphs improves the transparency of LLM  
150 reasoning (Da et al., 2024). However, detecting  
151 uncertainty is only the first step. The key challenge  
152 is mitigating or resolving it.

153 One approach to reducing ambiguity in LLMs  
154 is the generation of clarifying questions, where the  
155 model requests additional context to improve confi-  
156 dence in its response (Zhang and Choi, 2025). An-  
157 other strategy involves model ensembles, which  
158 aggregate outputs from multiple independently  
159 trained LLMs with different parameters and archi-  
160 tectures, reducing overall uncertainty (Liu et al.,  
161 2024). Additionally, embeddings from small lan-  
162 guage models can be used to distinguish between  
163 different types of uncertainty (Ahdritz et al., 2024).  
164 Methods based on CP have also been proposed to  
165 calibrate uncertainty in LLM-based planning sys-  
166 tems (Angelopoulos and Bates, 2022).

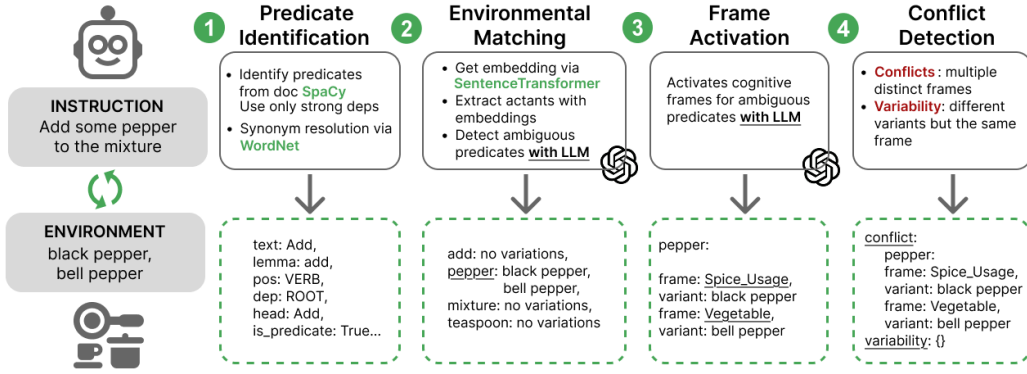


Figure 2: This figure illustrates the pipeline for detecting semantic ambiguity of the SVC method. The architecture decomposes the detection process into four main stages, each corresponding to a step in the cognitive mechanism of frame activation and conflict evaluation. The outputs of each stage of pipeline is illustrated using the example instruction: “Add some pepper to the mixture” and environment: black pepper, bell pepper. The system identifies potential ambiguity by tracing how the lexical unit pepper may evoke distinct semantic interpretations depending on context.

To address ambiguity resolving in robotic instruction-following, various approaches have been introduced in LLM-based planning. KnowNo Ren et al. (2023) is a framework that enables planners to assess and align uncertainty, helping them recognize when they lack confidence and need external input. This ensures statistically reliable task completion through CP. Building on this, Introspective Planning (IntroPlan) (Liang et al., 2024) integrates retrieval-augmented planning with CP, allowing models to proactively assess their confidence before taking action. By doing so, it reduces the number of user queries for task clarification. Another complementary approach is LAP (Jr. and Manocha, 2024), which includes the A-Feasibility metric. This metric combines scene context and model prompting to evaluate whether an action is both feasible and safe in environment. Ambiguities in large, shared spaces often arise from underspecified instructions that depend on implicit semantic features (e.g. cleanliness, fullness). To address this, Dogan et al. (2025) propose a model-agnostic approach leveraging iterative clarifications grounded in knowledge embeddings to infer missing attributes and improve object localization. A related line of work by Jiang et al. (2025) emphasizes that such vagueness frequently originates from referring expressions whose meaning is shaped by dialogue context and environmental factors.

In our work, we address a specific and underexplored source of ambiguity: semantic ambiguity in natural language instructions caused by lexical polysemy and underspecified argument roles. We treat ambiguity as a structural conflict between compet-

ing semantic frames activated by the same lexeme.

### 3 Background

In this section, we introduce key concepts used throughout the paper, **illustrative examples** are provided in Appendix A. Research in cognitive linguistics has shown that word meanings are shaped by underlying conceptual structures (Schwarze and Schepping, 1995; Bierwisch, 1983; Lehrer, 1990). Following frame semantics (Fillmore, 1985), we assume that understanding a lexeme requires reference to a conceptual **frame** specifying its semantic roles. Central to this view is the **lexical core** – the invariant meaning around which the lexeme’s interpretations are structured.

We define the **semantic valency** of a lexeme  $L$  as the set of independent conceptual variables  $X$  necessary for interpreting  $L$ ’s core meaning. Lexemes that require one or more such variables are **predicates**, and the expressions that realize them in context are **semantic actants** (Testelefs, 2001). Actants instantiate the predicate’s roles in specific situations, linking conceptual meaning with linguistic realization. For instance, in the instruction “Pour me some tea”, the predicate tea requires specification of a type (e.g., green or black). In a context where multiple instantiations of this actant are possible, such as green tea and black tea co-occurring in the environment, the valency slot remains undetermined. Our method identifies this as a semantic conflict: multiple actants (types of tea) compete within the same frame (Type\_of\_drink), revealing ambiguity rooted in underspecified preference (see example in Fig. 1. Our ambiguity detection module follows Yarowsky (1993), who observed that

ambiguity often arises when a predicate activates incompatible valency structures simultaneously.

## 4 Method

**Problem Formulation.** Formally, let  $L$  be a lexeme (e.g. a verb, noun or pronoun), and let  $X = \{x_1, x_2, \dots, x_n\}$  be the set of semantic valencies of  $L$ , where each  $x_i$  corresponds to an argument slot (actant) required for coherent interpretation. We define a **Semantic Valency Conflict** (SVC) as a case where multiple interpretations of  $L$  are simultaneously activated, and their corresponding valency profiles are incompatible with respect to the contextual environment  $\mathcal{E}$ . The environment  $\mathcal{E}$  is the set of objects, properties, and relations present in the external context where the interpretation occurs. It provides the grounding necessary to resolve or exacerbate potential conflicts between valency profiles.

Each interpretation  $I_k$  of  $L$  activates a cognitive frame  $\mathcal{F}_k$ , which implicitly defines a *valency profile*  $\mathcal{V}_k \subset \mathcal{D}$  – a set of expected participant roles and their semantic types. The domain  $\mathcal{D}$  represents the space of all conceptually valid participant configurations for predicates in natural language.

To determine whether two valency profiles  $\mathcal{V}_k$  and  $\mathcal{V}_m$  are semantically compatible in a given environment  $\mathcal{E}$ , we use the relation  $\sim_{\mathcal{E}}$ . This relation denotes approximate contextual equivalence:  $\mathcal{V}_k \sim_{\mathcal{E}} \mathcal{V}_m$  if their roles can be grounded to overlapping or compatible entities in  $\mathcal{E}$ , based on semantic similarity. In practice, this is operationalized using vector representations of candidate actants, their modifiers, and their types, matched against the entities and relations present in the environment.

If there exists  $k \neq m$  such that  $\mathcal{V}_k \not\sim_{\mathcal{E}} \mathcal{V}_m$ , then  $L$  is ambiguous. For a detailed illustration of the valency conflict detection process, see Fig. 2.

It is important to note that in our approach the term predicate is used in a broad sense, encompassing not only verbs but also nouns or even entire phrases that activate distinct cognitive frames.

**Overview of the SVC Method.** To address the problem of detecting ambiguity arising from the lexical polysemy and underspecified argument structures, we propose the **Semantic Valency Conflict** method. It is based on the assumption that ambiguity in a natural language instruction is determined by a conflict between different semantic valencies of the same lexeme (Long et al., 2022; Apresjan, 2000). In other words, ambiguity arises

when a single lexeme simultaneously activates multiple incompatible frames (cognitive scenarios)  $\mathcal{F}$ .

To implement the SVC method, we propose an architecture (see Fig. 2) where each stage addresses a specific task, from preprocessing to ambiguity detection. The overall architecture mirrors the cognitive process of frame activation and conflict resolution, and is implemented in four sequential stages. More details of each stage of SVC architecture and **algorithm** are provided in Appendix E.

The actants of each predicate are extracted and aligned with environmental constraints  $\mathcal{E}$  (i.e., the set of relevant objects, properties, and relations in the instruction’s context). The goal is to determine whether all participant roles required by the predicate are contextually supported. For each possible interpretation  $I_k$  of a predicate, the corresponding cognitive frame  $\mathcal{F}_k$  is activated. Frame induction is performed dynamically using LLMs, allowing the system to infer plausible frames and their valency structures without relying on static knowledge bases (Liang et al., 2024; Chaplot and Salakhutdinov, 2018). Each frame implicitly defines a valency profile  $\mathcal{V}_k$  over the domain  $\mathcal{D}$ .

At the final stage, the system evaluates the set of cognitive frames  $\{\mathcal{F}_k\}$  activated for each predicate, each associated with a valency profile  $\mathcal{V}_k$ . The method compares valency profiles  $\mathcal{V}_k$  for all competing interpretations. If there exist two interpretations  $k \neq m$  such that their valency profiles are semantically incompatible with respect to the environment, formally expressed as

$$\mathcal{V}_k \not\sim_{\mathcal{E}} \mathcal{V}_m, \quad (1)$$

the predicate is flagged as ambiguous due to **frame conflict**. This condition indicates that the predicate simultaneously activates multiple, mutually exclusive event schemas, reflecting true semantic ambiguity – for example, the predicate *pepper* activating frames  $\mathcal{F}_1 = \text{Spice\_Usage}$  and  $\mathcal{F}_2 = \text{Vegetable\_Cooking}$ .

Alternatively, if only one frame  $\mathcal{F}_k$  is active but multiple variants  $\{v_i\}$  exist within this frame,

$$|\{v_i\}| > 1, \quad (2)$$

the system interprets this as a **choice within the scenario** rather than a semantic conflict. Such variability represents different referents or attributes within the same event schema (e.g., variants *olive oil* and *sunflower oil* under the frame  $\text{Cooking\_Oil\_Ingredient}$ ). Although this does

not signify ambiguity at the frame level, it highlights distinctions that must be resolved for precise grounding or execution. Thus, the ambiguity detection mechanism differentiates between the cases providing a robust foundation for subsequent disambiguation steps.

## 5 Experimental Evaluation

To systematically assess the effectiveness and limitations of the proposed approach, we conduct a comprehensive experimental evaluation structured around the following research questions (RQs):

- **RQ1:** Can semantic conflict serve as a sufficient basis for ambiguity detection?
- **RQ2:** In which task types is semantic valency conflict most significant?
- **RQ3:** How do model size/type/architecture affect method performance?
- **RQ4:** Is simple prompting competitive?

**Datasets.** To evaluate SVC, we used two datasets: **IntroPlan Mobile Manipulation** (Liang et al., 2024) and **AmbiK** (Ivanova et al., 2025). Both contain robot instructions within a given environment and were used only for evaluation; no training or tuning was performed on them.

IntroPlan Mobile Manipulation dataset includes 600 short, single-step tasks with scene descriptions and object lists with the same distribution of different types of examples as in the KnowNo dataset (Ren et al., 2023). Although it contains ambiguous examples, it lacks precise ambiguity definitions and has an incomplete typology that overlaps with general linguistic phenomena.

AmbiK is a specialized dataset containing 2,000 annotated instructions covering three types of semantic ambiguity: user preferences, commonsense knowledge, and safety considerations. Each task has clarifying questions, answers, and execution plans for both ambiguous and unambiguous versions, and success markers for disambiguation.

**Metrics.** To evaluate the quality of ambiguity detection, we employ two binary metrics from Ivanova et al. (2025): **Help Rate (HR)** and **Correct Help Rate (CHR)**. **HR** measures how often the detection module determines that an instruction requires intervention – such as a clarifying question or a help request. The signal for this decision is the presence of either a semantic conflict between

activated frames or the variability of possible actant interpretations within a single frame. **CHR** assesses the appropriateness of the help request decision, taking into account the actual ambiguity type annotated in the dataset. For unambiguous instructions, the robot should not request help; conversely, it should request help when ambiguity is present, depending on the type of instruction.

To assess the alignment between system behavior and user intent, we additionally report **Intent Alignment (IA)**, which captures whether the method’s interpretation or suggested disambiguation correctly preserves the user’s intended goal. IA is computed by embedding the set of ground-truth user intents and the set of interpretations produced by the method into a shared semantic space, and assigning a positive score if at least one predicted variant exceeds a cosine similarity threshold with ground-truth intent.

Finally, as a main metric, we introduce **Contextual Help Quality (CHQ)** to measure the usefulness of contextual clarifications in supporting user intent. CHQ evaluates the combined effect of accurate ambiguity detection and intent preservation, and is defined as the harmonic mean of IA and CHR:

$$\text{CHQ} = \frac{2 \cdot \text{IA} \cdot \text{CHR}}{\text{IA} + \text{CHR} + \varepsilon}, \quad (3)$$

where  $\varepsilon$  is a small constant ( $\varepsilon = 10^{-6}$ ) added for numerical stability. For more information on the metrics, see Appendix B.

**Models and Baselines.** We conducted experiments using **GPT-3.5-Turbo** (OpenAI, 2023)<sup>1</sup>, **Mistral-7B-Instruct-v0.2** (MistralAI, 2023)<sup>2</sup>, **Gemma-2-9B-IT** (Google, 2024)<sup>3</sup>, and **Qwen1.5-7B-Chat** (QwenTeam, 2024)<sup>4</sup>. For brevity and clarity, we refer to these models simply as GPT-3.5, Mistral, Gemma, and Qwen, respectively. To compare the proposed SVC-based approach with a prompt-only baseline, we additionally employed the recently released **GPT-5-mini** model (OpenAI, 2025)<sup>5</sup>. For implementation details and prompts used see Appendix C.

We compare our method against **four baselines** previously proposed in the context of instruction

<sup>1</sup>Accessed via API: <https://platform.openai.com>

<sup>2</sup>Available at: <https://huggingface.co/mistralai/Mistral-7B-Instruct-v0.2>

<sup>3</sup>Available at: <https://huggingface.co/google/gemma-2-9b-it>

<sup>4</sup>Available at: <https://huggingface.co/Qwen/Qwen1.5-7B-Chat>

<sup>5</sup>Accessed via API: <https://platform.openai.com>

Table 1: Performance across instruction types. For each model and metric, the best result is shown in **bold**, while the second-best result is underlined. Higher values are better for CHR, IA, and CHQ. Additional results with HR see in Appendix 4

Model	Method	Instruction type											
		Unambiguous			Common Sense			Preferences			Safety		
		CHR↑	IA↑	CHQ↑	CHR↑	IA↑	CHQ↑	CHR↑	IA↑	CHQ↑	CHR↑	IA↑	CHQ↑
GPT-3.5	LoFreeCP	<u>0.77</u>	0.18	0.29	<u>0.80</u>	0.10	0.17	<u>0.15</u>	<u>0.10</u>	<u>0.12</u>	<u>0.24</u>	0.10	<u>0.14</u>
	LAP	<b>1.00</b>	<b>0.41</b>	<b>0.58</b>	<b>1.00</b>	0.20	<u>0.33</u>	0.00	<u>0.10</u>	0.00	0.00	0.20	0.00
	KnowNo	<b>1.00</b>	<u>0.35</u>	<u>0.51</u>	<b>1.00</b>	<u>0.21</u>	<b>0.34</b>	0.00	<u>0.10</u>	0.00	0.00	<u>0.22</u>	0.00
	SVC	0.51	<u>0.35</u>	0.41	0.51	<b>0.26</b>	<b>0.34</b>	<b>0.77</b>	<b>0.58</b>	<b>0.66</b>	<b>0.43</b>	<b>0.26</b>	<b>0.32</b>
Mistral	LoFreeCP	0.28	<b>0.69</b>	<b>0.39</b>	0.31	<u>0.40</u>	<b>0.34</b>	<u>0.73</u>	<u>0.52</u>	<u>0.60</u>	<b>0.62</b>	<b>0.50</b>	<b>0.55</b>
	LAP	<u>0.93</u>	0.12	<u>0.21</u>	<u>0.95</u>	0.05	0.09	0.06	0.04	0.04	0.04	0.02	0.02
	KnowNo	<b>1.00</b>	0.02	0.03	<b>1.00</b>	0.00	0.00	0.00	0.03	0.00	0.00	0.00	0.00
	SVC	0.51	<u>0.32</u>	<b>0.39</b>	0.23	<b>0.43</b>	<u>0.30</u>	<b>0.77</b>	<b>0.74</b>	<b>0.75</b>	<u>0.44</u>	<u>0.45</u>	<u>0.44</u>
Qwen	LoFreeCP	0.36	<b>0.24</b>	<u>0.28</u>	0.26	<b>0.28</b>	0.27	<u>0.62</u>	<u>0.32</u>	<u>0.42</u>	0.17	<u>0.23</u>	0.19
	LAP	<u>0.73</u>	<u>0.22</u>	<b>0.33</b>	<u>0.64</u>	0.21	<u>0.31</u>	0.28	<u>0.17</u>	<u>0.21</u>	<u>0.29</u>	<u>0.17</u>	<u>0.21</u>
	KnowNo	<b>1.00</b>	0.08	0.14	<b>1.00</b>	0.02	0.03	0.00	0.19	0.00	0.00	0.00	0.00
	SVC	0.54	0.17	0.25	0.61	<u>0.27</u>	<b>0.37</b>	<b>0.71</b>	<b>0.37</b>	<b>0.48</b>	<b>0.39</b>	<b>0.25</b>	<b>0.30</b>
Gemma	LoFreeCP	0.51	<b>0.41</b>	<b>0.45</b>	0.44	<b>0.20</b>	<u>0.27</u>	<u>0.20</u>	<u>0.24</u>	<u>0.21</u>	<b>0.34</b>	<u>0.09</u>	<u>0.14</u>
	LAP	0.66	0.12	0.20	<u>0.89</u>	0.04	0.07	0.00	0.10	0.00	0.13	0.05	0.07
	KnowNo	<b>1.00</b>	0.09	0.16	<b>1.00</b>	0.00	0.00	0.00	0.17	0.00	0.00	0.00	0.00
	SVC	<u>0.83</u>	<u>0.26</u>	<u>0.39</u>	0.84	<u>0.17</u>	<b>0.28</b>	<b>0.37</b>	<b>0.54</b>	<b>0.43</b>	0.21	<b>0.32</b>	<b>0.25</b>

426 following and embodied agents. **KnowNo** (Ren  
427 et al., 2023) is one of the first methods to apply  
428 conformal prediction for ambiguity detection in  
429 kitchen task scenarios with LLMs. The model  
430 is prompted to generate multiple candidate inter-  
431 pretations and then select the most appropri-  
432 ate one. **LAP** (Jr. and Manocha, 2024) extends  
433 KnowNo by incorporating affordance estimation.  
434 For each candidate response, it combines the model  
435 confidence with two affordance scores, Context-  
436 Based Affordance and Prompt-Based Affordance  
437 **LoFreeCP** (Su et al., 2024b) avoids logit access  
438 and instead applies uncertainty-based conformal  
439 prediction over multiple LLM generations. In addi-  
440 tion, we include a **Prompting** baseline, represent-  
441 ing a standard instruction-based ambiguity detec-  
442 tion setup without conformal prediction or external  
443 supervision. The LLM is prompted to classify an  
444 instruction as Ambiguous or Unambiguous given  
445 the environment and, if ambiguous, to produce a  
446 minimal set of alternative interpretations (see Ap-  
447 pendix C). We consolidate the **Binary** and **NoHelp**  
448 baselines from (Ren et al., 2023) into this unified  
449 **Prompting** baseline, as both rely on direct prompt-  
450 ing and differ only in whether uncertainty is explic-  
451 itly expressed: **Binary** labels its output as Certain  
452 or Uncertain and requests clarification accordingly,  
453 whereas **NoHelp** always outputs a single answer

without requesting clarification.

454  
455 **Results.** Our experimental evaluation demon-  
456 strates that the proposed Semantic Valency Conflict  
457 method provides a robust, semantically grounded  
458 mechanism for detecting ambiguity in embodied in-  
459 structions. Across both datasets, SVC competitive  
460 with many established baselines, with particularly  
461 strong gains on instruction types related to safety-  
462 critical scenarios and unambiguous instructions,  
463 where precise predicate grounding is essential.

464 **RQ1: Semantic conflict is sufficient for ambi-**  
465 **guity detection.** The strong performance of SVC,  
466 particularly on safety-based and contextually un-  
467 derspecified instructions, empirically validates that  
468 conflict between induced semantic frames is a suf-  
469 ficient and robust signal. The AmbiK dataset (Ta-  
470 ble 1) reveals SVC’s principal strength: identifying  
471 ambiguity arising from user preference. In this cat-  
472 egory, SVC achieves a CHR of 0.77 (vs. 0.15 for  
473 LoFreeCP on GPT and 0.73 on Mistral).

474 **RQ2: Semantic valency conflict is most signif-**  
475 **icant in underspecified instruction types.** It is  
476 most critical for Underspecification tasks. These  
477 tasks involve a predicate whose actants can be re-  
478 alized by multiple distinct entities in the environ-  
479 ment. The conflict emerges from the choice be-  
480 tween these realizations. Table 2 shows that SVC

Table 2: Results on the IntroPlan dataset across instruction types and models. The table reports CHR, IA, and CHQ for Prompting, LAP, LoFreeCP, and SVC. For each model and metric, the best result is shown in **bold**, while the second-best result is underlined. The KnowNo baseline was excluded from this table because preliminary evaluation showed it consistently produced degenerate outputs (HR = 0) by failing to identify ambiguous instructions.

Instruction Type	Prompting			LAP			LoFreeCP			SVC		
	CHR↑	IA↑	CHQ↑	CHR↑	IA↑	CHQ↑	CHR↑	IA↑	CHQ↑	CHR↑	IA↑	CHQ↑
<b>Mistral-7B-Instruct</b>												
creative multilabel	0.00	0.77	0.00	0.19	0.22	<u>0.20</u>	1.00	0.00	0.00	0.22	1.00	<b>0.36</b>
single-label	1.00	0.68	<b>0.81</b>	0.80	0.25	0.38	0.00	0.00	0.00	0.40	0.95	<u>0.56</u>
winograd	1.00	0.73	<b>0.84</b>	0.30	0.24	<u>0.27</u>	1.00	0.13	0.23	0.00	0.00	0.00
multi-label	0.00	0.73	0.00	0.57	0.42	<u>0.48</u>	1.00	0.00	0.00	0.38	0.88	<b>0.53</b>
unambiguous	1.00	0.84	<u>0.91</u>	0.92	0.13	0.23	0.00	0.00	0.00	0.98	1.00	<b>0.99</b>
spatial amb	0.00	0.95	0.00	0.00	0.01	0.00	1.00	0.00	0.00	0.13	1.00	<b>0.23</b>
unsafe	0.00	0.67	0.00	0.08	0.17	<u>0.11</u>	1.00	0.00	0.00	0.44	1.00	<b>0.61</b>
creative single-label	1.00	0.79	<b>0.88</b>	0.73	0.31	0.44	0.00	0.00	0.00	0.69	0.50	<u>0.58</u>
<b>Qwen1.5-7B-Chat</b>												
creative multi-label	0.09	0.58	0.16	0.21	0.31	0.25	1.00	0.17	<u>0.29</u>	0.23	0.84	<b>0.36</b>
single-label	0.17	0.82	0.28	0.64	0.29	<u>0.40</u>	0.00	0.00	0.00	0.59	0.51	<b>0.55</b>
winograd	1.00	0.81	<b>0.90</b>	0.47	0.22	<u>0.30</u>	1.00	0.06	0.11	0.00	0.00	0.00
multi-label	0.00	0.84	0.00	0.63	0.53	<b>0.58</b>	1.00	0.00	0.00	0.27	0.86	<u>0.41</u>
unambiguous	0.51	0.35	<u>0.42</u>	0.84	0.15	0.25	0.00	0.00	0.00	0.34	0.72	<b>0.46</b>
spatial amb	0.16	0.47	<u>0.24</u>	0.00	0.04	0.00	1.00	0.00	0.00	0.33	0.21	<b>0.26</b>
unsafe	0.00	0.68	0.00	0.12	0.20	<u>0.15</u>	1.00	0.00	0.00	0.64	0.73	<b>0.68</b>
creative single-label	0.17	0.82	<u>0.28</u>	0.62	0.38	<b>0.47</b>	0.00	0.00	0.00	0.72	0.13	0.22
<b>Gemma-2-9B-IT</b>												
creative multi-label	0.00	0.00	0.00	0.14	0.28	0.19	1.00	0.14	<u>0.25</u>	0.23	0.84	<b>0.36</b>
single-label	1.00	0.13	0.23	0.56	0.22	<u>0.32</u>	0.00	0.00	0.00	0.59	0.51	<b>0.55</b>
winograd	0.00	0.00	0.00	0.34	0.22	<b>0.27</b>	1.00	0.00	0.00	0.00	0.00	0.00
multi-label	0.00	0.31	0.00	0.49	0.62	<b>0.55</b>	1.00	0.12	0.21	0.27	0.86	<u>0.41</u>
unambiguous	1.00	0.04	0.08	0.89	0.15	<u>0.26</u>	0.00	0.00	0.00	0.34	0.72	<b>0.46</b>
spatial amb	0.00	0.00	0.00	0.00	0.02	0.00	1.00	0.00	0.00	0.33	0.21	<b>0.26</b>
unsafe	0.00	0.00	0.00	0.10	0.19	<u>0.13</u>	1.00	0.04	0.08	0.64	0.73	<b>0.68</b>
creative single-label	1.00	0.00	0.00	0.76	0.34	<b>0.47</b>	0.00	0.00	0.00	0.72	0.13	<u>0.22</u>

481 outperforms almost all baselines on unambiguous  
482 and safety-related (unsafe) instruction types across  
483 all evaluated models. In unambiguous settings,  
484 SVC achieves near-perfect CHR and CHQ scores  
485 (e.g., up to 0.99 on Mistral-7B), indicating reliable  
486 rejection of false ambiguity, whereas prompting-  
487 based and uncertainty-driven baselines frequently  
488 over-trigger ambiguity signals. In safety-critical  
489 tasks, SVC yields substantially higher CHQ and  
490 IA scores, reflecting more accurate identification  
491 of ambiguity, while baselines miss such cases or  
492 collapse to degenerate behavior (CHR = 0).

493 **RQ3: SVC generalizes across architectures.**

494 Performance is consistent across model families,  
495 but quality is influenced by instruction-tuning.  
496 SVC’s logit-free design makes it portable. As seen  
497 in Table 1, it works effectively with GPT, Mistral,

498 Qwen, and Gemma. However, the quality of the in-  
499 duced frames and consequently the IA scores tends  
500 to be higher with instruction-tuned models (Mistral-  
501 7B-Instruct and Gemma-2-9B-IT). These models  
502 produce more structured, deterministic frame de-  
503 scriptions, which improves the downstream match-  
504 ing and conflict analysis.

505 **RQ4: Simple prompting is less effective**

506 **than SVC.** While simple prompting underper-  
507 forms SVC on structural ambiguity metrics (e.g.,  
508 CHQ), preference-based results are comparable  
509 (see Fig. 3). This is expected, as preference judg-  
510 ments rely primarily on surface-level plausibility  
511 rather than explicit decomposition of semantic con-  
512 flicts. In such settings, strong models can approx-  
513 imate the correct preference without isolating the  
514 underlying semantic conflict, reducing the observ-

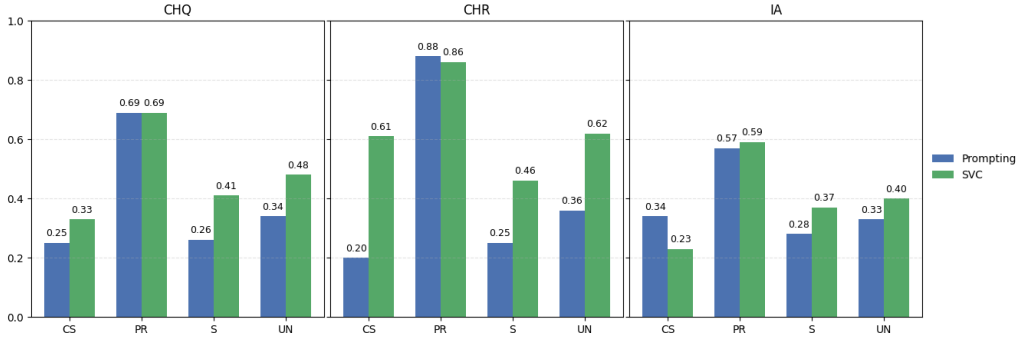


Figure 3: Comparison of CHQ performance between **Prompting** (prompt-only LLM baseline) and the proposed **SVC** method on the AmbiK dataset. Across all instruction types (Unambiguous (UN), Preference (PR), Common Sense (CS), and Safety (S)), SVC outperforms direct prompting in majority of cases, indicating improved robustness in ambiguity detection. The experiment is conducted using the modern **GPT-5-mini** model to reflect the current capabilities of large language models in handling complex ambiguities in natural language instructions.

Table 3: CHQ on a 150-sample AmbiK subset across all instruction types (Unambiguous (UN), Preference (PR), Common Sense (CS), and Safety (S)) using Mistral-7B. Best mean values are in **bold**. Abbreviations: P\_I – predicate identification, S\_E – synonym expansion, M\_H – modifier handling, E\_M – environment matching.

Ablation	CS	PR	S	UN
baseline	0.13±0.02	<b>0.60±0.03</b>	<b>0.41±0.04</b>	<b>0.43±0.02</b>
w/o P_I	0.13±0.02	0.58±0.04	0.39±0.05	<b>0.43±0.03</b>
w/o S_E	0.14±0.03	0.53±0.05	0.31±0.04	0.30±0.04
w/o M_H	0.06±0.02	0.58±0.04	0.31±0.03	0.26±0.03
w/o E_M	<b>0.30±0.04</b>	0.55±0.03	0.40±0.04	0.11±0.02

able gap between prompting and SVC.

## 6 Ablation Study

To further quantify the contribution of each module, we conducted systematic ablation studies to assess model sensitivity and component-level importance (see Table 3). The evaluation was performed on a 150-sample subset of the AmbiK dataset using the **Mistral-7B** model. To account for stochastic variability in LLM outputs, each ablation setting was evaluated over five independent runs with different random seeds, and the reported results are presented as the mean performance.

Pipeline exhibits strong robustness: removing most modules results in only moderate performance degradation. The most pronounced drop is observed when **environment matching** (E\_M) is removed. This leads to a sharp decline in Unambiguous (UN) performance (0.43 to 0.11), highlighting the critical role of explicit grounding in the environment for resolving referential and underspecification ambiguities. Interestingly, the same ablation yields an apparent improvement on Common Sense (CS) tasks (0.13 to 0.30). This suggests that CS ambiguities are primarily resolved at the semantic or conceptual level and do not require

explicit environment grounding. In such cases, environment matching may introduce spurious constraints, obscuring otherwise straightforward interpretations. **Synonym expansion** (S\_E) and **modifier handling** (M\_H) contribute consistent but secondary gains, particularly for CS and Safety tasks, supporting the importance of lexical and compositional refinement. In contrast, removing **predicate identification** (P\_I) results in minimal variation across all instruction types, indicating that the induced frame structure is largely stable once a predicate is identified.

These results confirm that SVC benefits from modular grounding mechanisms when task structure requires them, while remaining robust and in some cases more effective when such grounding is unnecessary.

## 7 Conclusion

In this work, we introduce **Semantic Valency Conflict**, a cognitively inspired method for detecting ambiguity in natural language instructions. The approach analyzes conflicts between semantic valency profiles generated by contextually activated cognitive frames. Method outperforms baselines on unambiguous and safety-related instruction types across all evaluated models, effectively identifies subtle forms of ambiguity. SVC requires neither access to model logits nor reliance on static knowledge bases, ensuring compatibility with both white-box and black-box systems while maintaining adaptability to novel environments and tasks. Moreover, SVC’s structured output provides interpretable diagnostics of ambiguity sources, facilitating effective downstream clarification and resolution mechanisms.

## 8 Limitations

Despite its effectiveness in detecting semantic valency conflicts, the current implementation of the SVC method is subject to several limitations.

**Environment completeness dependency.** The accuracy of ambiguity detection is highly contingent on the completeness of the environment model. If critical entities, properties, or relations are missing from  $\mathcal{E}$ , the system may fail to identify conflicts. Thus, noisy environment representation can undermine the method’s reliability.

**Restriction to prototypical valencies.** Our approach focuses on core, prototypical valency structures, excluding non-canonical and metaphorical usages. As a result, polysemous extensions – such as metaphorical interpretations or constructions with optional actants – are not captured by the current model and may lead to misclassification or under-detection of ambiguity. For example, the instruction “Kill the heat” is a metaphorical expression meaning turn off the stove or reduce the flame, but the model may misinterpret it due to the lack of a literal, manipulable object. Similarly, in “Season the soup”, the core valency structure (verb “season” + object “soup”) omits the specific substance to be added (e.g., salt, pepper, herbs), which is crucial for task execution but remains implicit and therefore may not be flagged as ambiguous.

**Syntactic and pragmatic ambiguity.** SVC’s performance collapses on Winograd tasks (CHR=0). These tasks hinge on coreference resolution and syntactic parsing. Since SVC operates on predicate-argument semantics and lacks a deep syntactic or discourse model, it cannot detect ambiguities where the semantic roles are stable but their referents are not.

**LLM dependency.** LLMs are used in multiple stages of the pipeline enhances the adaptability of the system compared to reliance on static knowledge bases, it introduces a trade-off between flexibility and precision. Specifically, LLMs are employed during frame activation to propose potential interpretations for ambiguous actions, and as a fallback mechanism to assess the ambiguity of actants whose reference to the environment remains uncertain. However, the primary actant extraction and entity alignment rely on symbolic parsing and embedding-based similarity, thereby reducing both computational load and epistemic dependency on

the LLM. This hybrid design strikes a balance between the adaptability offered by LLMs and the robustness of structured, similarity-driven matching. Nevertheless, the overall reliability of ambiguity detection still partially depends on the quality of LLM-generated outputs at these two stages, which may introduce noise or contextually inappropriate variants.

**Single-turn instruction processing.** The current system operates on isolated, single-turn instructions paired with a static environment description. It does not support multi-turn dialogue or interactive clarification.

## 9 Ethical Considerations

Our approach uses LLMs in generation mode. Even when using a prompt that limits the output format, the model may still generate inappropriate and/or offensive content. In addition, LLMs are prone to hallucinations and can produce unexpected results. Therefore, giving them control over potentially harmful machines, such as robots, and testing such machines should be done in a regulated manner in a specially designated area with limited access for the people involved in the experiments. It is also possible to deliberately instruct a robot to execute harmful actions.

## References

- Gustaf Ahndritz, Tian Qin, Nikhil Vyas, Boaz Barak, and Benjamin L. Edelman. 2024. [Distinguishing the knowable from the unknowable with language models](#). *Computing Research Repository*, arXiv:2402.03563. Version 2.
- Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, Chuyuan Fu, Keerthana Gopalakrishnan, Karol Hausman, and 1 others. 2022. Do as i can, not as i say: Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691*.
- Anastasios N. Angelopoulos and Stephen Bates. 2022. [A gentle introduction to conformal prediction and distribution-free uncertainty quantification](#). *Computing Research Repository*, arXiv:2107.07511.
- Jurij D. Apresjan. 2000. Regular polysemy and lexical functions. In Igor Mel’čuk and Leo Wanner, editors, *Systematic Lexicography*, pages 119–132. Oxford University Press, Oxford.
- Manfred Bierwisch. 1983. Semantische und conceptuelle repräsentation lexikalischer einheiten. In *Untersuchungen zur Semantik*. Akademie Verlag.

673	Devendra Singh Chaplot and Ruslan Salakhutdinov.	Wenlong Huang, Pieter Abbeel, Deepak Pathak, and	727
674	2018. <a href="#">Knowledge-based word sense disambiguation</a>	Igor Mordatch. 2022. Language models as zero-shot	728
675	<a href="#">using topic models</a> . <i>Computing Research Repository</i> ,	planners: Extracting actionable knowledge for em-	729
676	arXiv:1801.01900.	bodied agents. In <i>International conference on ma-</i>	730
677	Anfisa A Chuganskaya, Alexey K Kovalev, and Alek-	<i>chine learning</i> , pages 9118–9147. PMLR.	731
678	sandr Panov. 2023. The problem of concept learning	Anastasia Ivanova, Bakaeva Eva, Zoya Volovikova,	732
679	and goals of reasoning in large language models. In	Alexey Kovalev, and Aleksandr Panov. 2025. <a href="#">AmbiK:</a>	733
680	<i>International Conference on Hybrid Artificial Intelli-</i>	<a href="#">Dataset of ambiguous tasks in kitchen environment</a> .	734
681	<i>gence Systems</i> , pages 661–672. Springer.	In <i>Proceedings of the 63rd Annual Meeting of the</i>	735
682	Longchao Da, Tiejun Chen, Lu Cheng, and Hua Wei.	<i>Association for Computational Linguistics (Volume 1:</i>	736
683	2024. <a href="#">Llm uncertainty quantification through di-</a>	<i>Long Papers)</i> , pages 33216–33241, Vienna, Austria.	737
684	<a href="#">rectional entailment graph and claim level response</a>	Association for Computational Linguistics.	738
685	<a href="#">augmentation</a> . <i>Computing Research Repository</i> ,	Chenxi Jiang, Chuha Zhou, and Jianfei Yang. 2025.	739
686	arXiv:2407.00994. Version 2.	<a href="#">Rei-bench: Can embodied agents understand vague</a>	740
687	Fethiye Irmak Dogan, Maithili Patel, Weiyu Liu,	<a href="#">human instructions in task planning?</a> <i>Preprint</i> ,	741
688	Iolanda Leite, and Sonia Chernova. 2025. <a href="#">A model-</a>	arXiv:2505.10872.	742
689	<a href="#">agnostic approach for semantically driven disam-</a>	James F. Mullen Jr. and Dinesh Manocha. 2024. <a href="#">Lap, us-</a>	743
690	<a href="#">biguation in human-robot interaction</a> . <i>Preprint</i> ,	<a href="#">ing action feasibility for improved uncertainty align-</a>	744
691	arXiv:2409.17004.	<a href="#">ment of large language model planners</a> . <i>Computing</i>	745
692	Charles J. Fillmore. 1985. Frames and the semantics	<i>Research Repository</i> , arXiv:2403.13198.	746
693	of understanding. <i>Quaderni di Semantica</i> , 6(2):222–	Aleksei Konstantinovich Kovalev and Aleksandr Igo-	747
694	254.	vich Panov. 2022. Application of pretrained large	748
695	Roya Firoozi, Johnathan Tucker, Stephen Tian,	language models in embodied artificial intelligence.	749
696	Anirudha Majumdar, Jiankai Sun, Weiyu Liu, Yuke	In <i>Doklady Mathematics</i> , volume 106, pages S85–	750
697	Zhu, Shuran Song, Ashish Kapoor, Karol Hausman,	S90. Springer.	751
698	Brian Ichter, Danny Driess, Jiajun Wu, Cewu Lu, and	Adrienne Lehrer. 1990. Polysemy, conventionality and	752
699	Mac Schwager. 2023. <a href="#">Foundation models in robotics:</a>	the structure of the lexicon. <i>Cognitive Linguistics</i> .	753
700	<a href="#">Applications, challenges, and the future</a> . <i>Computing</i>	Kaiqu Liang, Zixu Zhang, and Jaime Fernández Fisac.	754
701	<i>Research Repository</i> , arXiv:2312.07843.	2024. Introspective planning: Aligning robots’ uncer-	755
702	Google. 2024. Gemma: Open models based on gemini	tainty with inherent task ambiguity. In <i>Proceedings</i>	756
703	research. Technical report.	<i>of the 38th Conference on Neural Information Pro-</i>	757
704	Danil Grigorev, Alexey Kovalev, and Aleksandr Panov.	<i>cessing Systems (NeurIPS)</i> .	758
705	2025. <a href="#">Verifyllm: Llm-based pre-execution task plan</a>	Linyu Liu, Yu Pan, Xiao Cheng Li, and Guanting Chen.	759
706	<a href="#">verification for robots</a> . In <i>2025 IEEE/RSJ Interna-</i>	2024. <a href="#">Uncertainty estimation and quantification for</a>	760
707	<i>tional Conference on Intelligent Robots and Systems</i>	<a href="#">llms: A simple supervised approach</a> . <i>Computing</i>	761
708	<i>(IROS)</i> .	<i>Research Repository</i> , arXiv:2404.15993.	762
709	Danil S Grigorev, Alexey K Kovalev, and Aleksandr I	Chen Long, Laura Kallmeyer, and Rainer Osswald.	763
710	Panov. 2024. Common sense plan verification with	2022. <a href="#">A frame-based model of inherent polysemy,</a>	764
711	large language models. In <i>International Conference</i>	<a href="#">copredication and argument coercion</a> . In <i>Proceed-</i>	765
712	<i>on Hybrid Artificial Intelligence Systems</i> , pages 224–	<i>ings of the Workshop on Cognitive Aspects of the</i>	766
713	236. Springer.	<i>Lexicon</i> , pages 58–67, Taipei, Taiwan. Association	767
714	Rishi Hazra, Pedro Zuidberg Dos Martires, and Luc De	for Computational Linguistics.	768
715	Raedt. 2024. Saycanpay: Heuristic planning with	MistralAI. 2023. <a href="#">Mistral-7b-instruct-v0.2.</a>	769
716	large language models using learnable domain knowl-	<a href="https://huggingface.co/mistralai/">https://huggingface.co/mistralai/</a>	770
717	edge. In <i>Proceedings of the AAI Conference on Arti-</i>	<a href="#">Mistral-7B-Instruct-v0.2</a> . Accessed: 2025-08-	771
718	<i>ficial Intelligence</i> , volume 38, pages 20123–20133.	01.	772
719	Juyeon Heo, Miao Xiong, Christina Heinze-Deml, and	Anatoly Onishchenko, Alexey Kovalev, and Aleksandr	773
720	Jaya Narain. 2025. Do llms estimate uncertainty well	Panov. 2025. Lookplangraph: Embodied instruction	774
721	in instruction-following? In <i>Proceedings of the In-</i>	following method with vlm graph augmentation. In	775
722	<i>ternational Conference on Learning Representations</i>	<i>Workshop on Reasoning and Planning for Large Lan-</i>	776
723	<i>(ICLR)</i> .	<i>guage Models</i> .	777
724	Matthew Honnibal, Ines Montani, Sofie Van Lan-	OpenAI. 2023. <a href="#">Gpt-3.5-turbo (august 16 version)</a> .	778
725	degheem, and Adriane Boyd. 2020. <a href="#">spacy: Industrial-</a>	<a href="https://openai.com">https://openai.com</a> . Accessed: 2024-08-16.	779
726	<a href="#">strength natural language processing in python</a> .	OpenAI. 2025. Gpt-5 technical report. Accessed via	780
		OpenAI API.	781

782	QwenTeam. 2024. Qwen1.5-7b-chat. <a href="https://huggingface.co/Qwen/Qwen1.5-7B-Chat">https://huggingface.co/Qwen/Qwen1.5-7B-Chat</a> . Accessed: 2025-08-01.	837
783		838
784		839
785	Nils Reimers and Iryna Gurevych. 2019. <a href="#">Sentence-bert: Sentence embeddings using siamese bert-networks</a> . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing</i> . Association for Computational Linguistics.	840
786		841
787		842
788		843
789		
790	Allen Z. Ren, Anushri Dixit, Alexandra Bodrova, Sumeet Singh, Stephen Tu, Noah Brown, Peng Xu, Leila Takayama, Fei Xia, Jake Varley, Zhenjia Xu, Dorsa Sadigh, Andy Zeng, and Anirudha Majumdar. 2023. Robots that ask for help: Uncertainty alignment for large language model planners. In <i>Proceedings of the Conference on Robot Learning (CoRL)</i> .	
791		
792		
793		
794		
795		
796		
797		
798	Christina Sarkisyan, Alexandr Korchemnyi, Alexey K Kovalev, and Aleksandr I Panov. 2023. Evaluation of pretrained large language models in embodied planning tasks. In <i>International Conference on Artificial General Intelligence</i> , pages 222–232. Springer.	
799		
800		
801		
802		
803	Christa Schwarze and Marie-Theres Schepping. 1995. Polysemy in a two-level-semantics. In <i>Current Issues in Linguistic Theory: Lexical Knowledge in the Organization of Language</i> . John Benjamins Publishing Company.	
804		
805		
806		
807		
808	Ola Shorinwa, Zhiting Mei, Justin Lidard, Allen Z. Ren, and Anirudha Majumdar. 2024. <a href="#">A survey on uncertainty quantification of large language models: Taxonomy, open research challenges, and future directions</a> . <i>Computing Research Repository</i> , arXiv:2412.05563.	
809		
810		
811		
812		
813		
814	Anmol Singhal, Chirag Jain, Preethu Rose Anish, Arkajyoti Chakraborty, and Smita Ghaisas. 2024. <a href="#">Generating clarification questions for disambiguating contracts</a> . <i>Computing Research Repository</i> , arXiv:2403.08053.	
815		
816		
817		
818		
819	Jiayuan Su, Jing Luo, Hongwei Wang, and Lu Cheng. 2024a. <a href="#">Api is enough: Conformal prediction for large language models without logit-access</a> . <i>Preprint</i> , arXiv:2403.01216.	
820		
821		
822		
823	Jiayuan Su, Jing Luo, Hongwei Wang, and Lu Cheng. 2024b. <a href="#">Api is enough: Conformal prediction for large language models without logit-access</a> . <i>Preprint</i> , arXiv:2403.01216.	
824		
825		
826		
827	Yakov G. Testeleets. 2001. <i>Introduction to General Syntax</i> . Russian State University for the Humanities, Moscow.	
828		
829		
830	Yifei Wang, Yu Sheng, Linjing Li, and Daniel Zeng. 2025. <a href="#">Uncertainty unveiled: Can exposure to more in-context examples mitigate uncertainty for large language models?</a> <i>Preprint</i> , arXiv:2505.21003.	
831		
832		
833		
834	David Yarowsky. 1993. One sense per collocation. In <i>Proceedings of the Workshop on Human Language Technology</i> .	
835		
836		
		844
		845
		846
		847
		848
	Hangtao Zhang, Chenyu Zhu, Xianlong Wang, Ziqi Zhou, Changgan Yin, Minghui Li, Lulu Xue, Yichen Wang, Shengshan Hu, Aishan Liu, Peijin Guo, and Leo Yu Zhang. 2025. <a href="#">Badrobot: Jailbreaking embodied LLMs in the physical world</a> . In <i>The Thirteenth International Conference on Learning Representations</i> .	
	Michael J. Q. Zhang and Eunsol Choi. 2025. Clarify when necessary: Resolving ambiguity through interaction with lms. In <i>Findings of the 2025 North American Chapter of the Association for Computational Linguistics (NAACL)</i> .	

## 849 A Appendix A – Glossary of Terms

- 850 • **Frame (Semantic Frame / Cognitive Frame)** A structured conceptual schema that organizes  
851 knowledge about a specific type of situation, event, or object, including the roles and relationships  
852 involved. Frames enable humans to interpret word meanings by invoking entire cognitive scenarios.  
853 *Example:* The sentence “The seller gave the customer some apples for five dollars” activates a  
854 COMMERCE frame involving roles such as SELLER, BUYER, GOODS, MONEY, and TRANSFER  
855 ACTION.
- 856 • **Cognitive Frame** Emphasizes the mental and experiential nature of a frame as a structure in the  
857 speaker’s conceptual system.  
858 *Example:* Interpreting the word “pepper” involves activating a mental model of either a cooking  
859 spice or a vegetable, depending on context.
- 860 • **Semantic Frame** Highlights the lexical-semantic dimension of frames, focusing on how word  
861 meanings and their valency structures relate to structured scenarios.  
862 *Example:* “Black pepper” activates the SPICE\_USAGE frame; “bell pepper” activates the VEG-  
863 ETABLE\_COOKING frame.
- 864 • **Cognitive Scenario** A dynamic instantiation of a frame, often used synonymously, placing emphasis  
865 on the sequence of events and interactions among participants.  
866 *Example:* A request like “pour me some tea” evokes a scenario involving serving, a recipient, and a  
867 beverage.
- 868 • **Predicate** A lexeme that requires one or more participants to complete its meaning, forming the  
869 semantic core around which arguments are organized.  
870 *Example:* In “pour me some tea,” the verb *pour* functions as a predicate that expects arguments like  
871 a recipient and a substance.
- 872 • **Lexical Core** The central, invariant meaning of a polysemous lexeme, abstracted from contextual  
873 variation. It provides the conceptual foundation from which semantic valencies emerge.  
874 *Example:* The noun *bank* has different senses (“financial institution” vs. “riverbank”), but in the  
875 financial sense, the lexical core refers to an institution that holds and manages money, regardless of  
876 specific services or context.
- 877 • **Semantic Valency** The set of argument slots that a lexeme (typically a predicate) conceptually  
878 requires in order to be fully interpreted. These slots correspond to independent semantic variables.  
879 *Example:* In the utterance “pour me some tea,” the predicate *tea* has a unary valency slot for *type of*  
880 *tea*, filled by either *green* or *black*.
- 881 • **Valency** An expression in a sentence that fills one of the semantic valency slots of a predicate.  
882 *Example:* *Green* and *black* are arguments for the valency slot in *tea* denoting type.
- 883 • **Semantic Actant** A linguistic expression that realizes one of the predicate’s semantic roles in a  
884 sentence. Actants serve as syntactic and semantic instantiations of frame roles.  
885 *Example:* In “The robot gave the customer tea,” the expressions *robot*, *customer*, and *tea* are actants  
886 corresponding to AGENT, RECIPIENT, and THEME roles, respectively.

## 887 B Appendix B – Metrics Details

888 **Intent Alignment (IA)** evaluates whether the clarification request issued by the model is *useful*, i.e.,  
889 whether it exposes alternative interpretations that are aligned with the user’s underlying intent annotation.  
890 IA is a binary proxy metric and does not aim to measure full intent recovery, but rather whether at least  
891 one generated clarification aligns with the annotated user intent.

892 IA is computed only for instructions where the model requested help (HR = 1). For all other cases, IA  
893 is undefined and excluded from aggregation.

---

**Algorithm 1: Intent Alignment (IA)**

---

**Input:** User intent annotation  $I$ , generated variants  $V$ , embedding model, threshold  $\tau$

**Output:**  $IA \in \{0, 1\}$

**if**  $I = \emptyset$  **then**

  | **return** NaN

**else**

  | Encode  $I$  and  $V$  using the embedding model;

  | Compute cosine similarities  $S_{ij}$  for all  $i \in I, j \in V$ ;

  |  $s_{\max} \leftarrow \max_{i,j} S_{ij}$ ;

  | **if**  $s_{\max} \geq \tau$  **then**

    | **return** 1

  | **else**

    | **return** 0

---

**Help Rate (HR)** HR measures how often the detection module determines that an instruction requires intervention – such as a clarifying question or a help request. The signal for this decision is the presence of either a *semantic conflict* between activated frames or *variability of possible actant interpretations* within a single frame. Formally,

$$HR = \begin{cases} 1, & \text{if detected ambiguous,} \\ 0, & \text{otherwise.} \end{cases} \quad (4)$$

**Correct Help Rate (CHR)** CHR assesses the appropriateness of the help request decision, taking into account the actual ambiguity type annotated in the dataset.

For unambiguous instructions, the robot should not request help; conversely, it should request help when ambiguity is present. Thus, CHR measures the *selectivity and accuracy* of the detection module. It is defined as follows depending on the ambiguity type.

The **AmbiK dataset** contains the following instruction types: *user preferences, common sense knowledge, safety* and *unambiguous tasks*.

- For tasks of type *user preferences* and *safety*, the model is considered correct if it requested help:

$$CHR = \begin{cases} 1, & \text{if } HR = 1, \\ 0, & \text{otherwise.} \end{cases} \quad (5)$$

- For other types (*common sense knowledge*, as well as *unambiguous tasks*), the model should succeed without requesting help:

$$CHR = \begin{cases} 1, & \text{if } HR = 0, \\ 0, & \text{otherwise.} \end{cases} \quad (6)$$

The **IntroPlan dataset** contains the following instruction types: *creative\_multilabel\_task, singlelabel\_task, winograd\_task, multilabel\_task, unambiguous\_task, spatial\_ambiguous\_task, unsafe\_task*, and *creative\_singlelabel\_task*.

- We consider as ambiguous *creative\_multilabel\_task, winograd\_task, multilabel\_task, unsafe\_task* and *spatial\_ambiguous\_task*, instructions where clarifications are generally expected

$$CHR = \begin{cases} 1, & \text{if } HR = 1, \\ 0, & \text{otherwise.} \end{cases} \quad (7)$$

- For other types (*unambiguous\_task*, *singlelabel\_task*, *creative\_singlelabel\_task*), the model should succeed without requesting help:

$$\text{CHR} = \begin{cases} 1, & \text{if HR} = 0, \\ 0, & \text{otherwise.} \end{cases} \quad (8)$$

This differentiated definition of CHR allows the evaluation to account for the specific nature of various ambiguity types: clarifications are necessary for user preference ambiguities, while for tasks relying on common knowledge or obvious actions, a help request may be redundant.

## C Appendix C – Implementation Experimental Setup

### C.1 Implementation Details

**Instruction Preprocessing.** To ensure consistency and facilitate downstream analysis, each instruction was preprocessed using a transformer-based spaCy (Honnibal et al., 2020) pipeline (en\_core\_web\_trf). This step includes tokenization, lemmatization, part-of-speech tagging, and dependency parsing. For each token in the instruction, we extracted linguistic features such as surface form, lemma, part-of-speech tag, syntactic dependency label, syntactic head, and stopword status.

**Semantic Matching via Embeddings.** To align mentioned arguments (actants) in the instruction with entities in the environment, we used semantic similarity in embedding space. Specifically, we employed the Sentence-BERT model (Reimers and Gurevych, 2019) (paraphrase-mpnet-base-v2) to encode both the textual mentions in the instruction and the object names in the environment. The resulting embeddings were compared using cosine similarity to identify the most relevant environmental entities for each mentioned actant.

**Language Model Configuration.** All experiments were conducted using GPT-3.5-Turbo, GPT-5-mini (via the OpenAI API), mistralai/Mistral-7B-Instruct-v0.2, Qwen/Qwen1.5-7B-Chat and google/gemma-2-9b-it. For GPT-3.5-Turbo and GPT-5-mini, we used a temperature of 0.3 and a maximum token limit of 400. The model was prompted using a consistent interaction format: a system message (“*You are a helpful assistant.*”) followed by a single user message containing the input prompt. For other models we used ‘do\_sample=True’ and ‘top\_p=0.95’.

### C.2 Prompt for Baseline Prompting

```
You are an ambiguity detection system.
You MUST return a JSON object with EXACTLY this structure:
"label": "Ambiguous" | "Unambiguous",
"ambiguities": [
"phrase": string,
"candidates": [string, string]
]
Rules:
- If there is ANY uncertainty, label = "Ambiguous"
- If label is "Unambiguous", ambiguities MUST be an empty list
- Use ONLY the given environment
- Do NOT add any extra keys
- Do NOT include explanations
- Return ONLY JSON
- Candidates must be short alternative interpretations, not explanations
- Do NOT describe feasibility, difficulty, or reasoning
- Candidates must be minimal noun-phrase interpretations only
```

Lexeme: lemma

Variants: variants\_str

Environment: environment\_str

Task:

Analyze each variant of the lexeme 'lemma' in the given environment and identify the cognitive frame it activates.

For each variant:

- Identify the cognitive frame it activates.
- List the expected semantic roles (valency) for that frame.

Important:

If the only differences between variants are minor attributes like color, material, or size, and these do not affect the core function or role of the object in the environment, then all variants should activate the same cognitive frame and share the same valency.

Only distinguish frames if the environment suggests clearly different roles or functions for the variants.

Each output block must include:

1. The variant
2. The name of the activated frame (e.g., "Drinkware\_Container")
3. A list of semantic roles (e.g., container, liquid, function)

**\*\*CRUCIAL:\*\***

- Output one variant per block, starting with a dash and a space ('- variant\_name:').
- Never list multiple variants in the same block or line (no commas in variant names).
- Separate blocks clearly with a blank line.

Example:

Lexeme: soda

Variants: coke, pepsi, orange soda

Environment: coke, pepsi, orange soda, bottled water

- coke:

frame: Flavored\_Soft\_Drink

valency: beverage, brand, carbonation

- pepsi:

frame: Flavored\_Soft\_Drink

valency: beverage, brand, carbonation

- orange soda:

frame: Citrus\_Drink

valency: beverage, fruit\_flavor, carbonation

## C.4 Prompt for Actant Extraction and Environment Matching

You are a semantic disambiguation assistant.

Here is a list of actants extracted from an instruction:

Actants: actant\_text

And a list of objects present in the current environment:

Environment: environment\_text

Your tasks:

1. For each actant, decide if it is ambiguous in the current environment.
2. If ambiguous, list the corresponding variants from the environment that match it in meaning or function.
3. If unambiguous, write "no variations".

Return your answer as a JSON object with actants as keys.

Example:

```
"cola": ["pepsi", "coke"],
"bottled water": "no variations"
```

## D Appendix D – Additional Evaluation Results

Table 4: HR and CHR scores across instruction types on Ambik dataset.

Model	Method	Instruction type							
		Unambiguous		Common Sense		Preferences		Safety	
		HR	CHR↑	HR	CHR↑	HR	CHR↑	HR	CHR↑
GPT	LoFreeCP	0.23	0.77	0.20	0.80	0.15	0.15	0.24	0.24
	LAP	0.00	1.00	0.00	1.00	0.00	0.00	0.00	0.00
	KnowNo	0.00	1.00	0.00	1.00	0.00	0.00	0.00	0.00
	SVC	0.49	0.51	0.48	0.51	0.77	0.77	0.43	0.43
Mistral	LoFreeCP	0.72	0.28	0.69	0.31	0.73	0.73	0.62	0.62
	LAP	0.07	0.93	0.05	0.95	0.06	0.06	0.04	0.04
	KnowNo	0.00	1.00	0.00	1.00	0.00	0.00	0.00	0.00
	SVC	0.49	0.51	0.77	0.23	0.77	0.77	0.44	0.44
Qwen	LoFreeCP	0.64	0.36	0.74	0.26	0.62	0.62	0.17	0.17
	LAP	0.27	0.73	0.36	0.64	0.28	0.28	0.29	0.29
	KnowNo	0.00	1.00	0.00	1.00	0.00	0.00	0.00	0.00
	SVC	0.46	0.54	0.39	0.61	0.71	0.71	0.39	0.39
Gemma	LoFreeCP	0.49	0.51	0.56	0.44	0.20	0.20	0.34	0.34
	LAP	0.34	0.66	0.11	0.89	0.00	0.00	0.13	0.13
	KnowNo	0.00	1.00	0.00	1.00	0.00	0.00	0.00	0.00
	SVC	0.17	0.83	0.16	0.84	0.37	0.37	0.21	0.21

## E Appendix E – Details of Ambiguity Detection of the SVC Method

In order to clarify the stepwise functioning of the SVC method, we now examine its implementation in detail:

---

**Algorithm 2:** Ambiguity Detection via Semantic Valency Conflicts

---

**Input:** Instruction  $I$ , environment objects  $\mathcal{E}$

**Output:** Ambiguity type (clear/variable/conflict), semantic variants

**1. Linguistic Analysis;**

$\mathcal{P} \leftarrow$  extract predicates from  $I$ ;

$\mathcal{A} \leftarrow$  extract actants for each  $p \in \mathcal{P}$ ;

**2. Multi-Stage Referent Matching;**

**foreach** actant  $a \in \mathcal{A}$  **do**

**if** exact match between  $a$  and  $e \in \mathcal{E}$  **then**

$R_a \leftarrow \{e\}$  // Single referent

**end**

**else**

$R_a \leftarrow \{e \in \mathcal{E} \mid \text{similarity}(a, e) > \tau\}$  // Possible referents

**if**  $|R_a| > 1$  **then**

            mark  $a$  as ambiguous

**end**

**end**

**end**

**3. Semantic Frame Activation;**

**if** any actant is ambiguous **then**

$\mathcal{V} \leftarrow$  generate semantic variants via LLM;

$\mathcal{F} \leftarrow$  activate frames for each variant;

**end**

**4. Conflict Analysis;**

**foreach** predicate  $p \in \mathcal{P}$  **do**

**if** multiple incompatible frames for  $p$  **then**

        report **conflict** // Ex: "cut"  $\rightarrow$  slice/chop

**else**

**if** multiple compatible frames for  $p$  **then**

            report **variability** // Ex: "bread"  $\rightarrow$  toast/sandwich

**else**

            report **clear**

**end**

**end**

**end**

**5. Intent Alignment;**

**if** user intent  $U$  provided **then**

    compute alignment between  $U$  and  $\mathcal{V}$ ;

**end**

**return** ambiguity type, semantic variants  $\mathcal{V}$ ;

---